

# NOT FOR PUBLIC RELEASE

1 Michael S.A. Graziano

2 *We Are Machines That*  
3 *Claim to Be Conscious*

4 **Abstract:** *The attention schema theory explains how a biological,*  
5 *information processing machine can claim to have consciousness, and*  
6 *how, by introspection (by assessing its internal data), it cannot deter-*  
7 *mine that it is a machine whose claims are based on computations.*  
8 *The theory directly addresses Chalmers' meta-problem of conscious-*  
9 *ness, the problem of why we think we have a difficult-to-explain*  
10 *consciousness in the first place.*

11 **1. Introduction**

12 Neuroscience has taught us that the brain is an information processing  
13 device. In the perspective that I take, and the theory I have suggested  
14 — the attention schema theory (AST) — we are information pro-  
15 cessing machines that, among other actions, make claims about our-  
16 selves (e.g. Graziano and Kastner, 2011; Graziano, 2013; Webb and  
17 Graziano, 2015). We claim to have something inside us, subjective  
18 experience, that is fundamentally non-physical. Logically, the brain  
19 cannot put out a claim unless it contains the information on which the  
20 claim is based. In my research, therefore, I have focused on the  
21 information set on which the claim of subjective experience is based.  
22 What cognitive purpose does it serve? What brain regions might be  
23 involved in constructing it? How is the machine engineered such that  
24 it makes that claim?

25 I would like to clarify at the outset what I mean by a non-physical  
26 property. If a person looks at a red apple, she not only processes  
27 information about the colour, but also claims to have a subjective  
28 experience of red — the ‘what it feels like’ component. One cannot  
29 push on subjective experience and measure a reaction force, scratch it

Correspondence:  
Email: [graziano@princeton.edu](mailto:graziano@princeton.edu)

1 and measure its hardness, or put it on a scale and measure its weight.  
2 It does not exist on those physical dimensions. In the sense of its  
3 physical non-measurability, subjective experience is non-physical, or  
4 even metaphysical in the strict sense of being above or outside the  
5 physical. This ethereal nature of subjective experience is precisely  
6 why it has been so difficult to understand.

7 But, objectively speaking, the phenomenon that faces us is much  
8 simpler. A brain-controlled agent constructs a self-description and on  
9 that basis makes claims about itself. There is no rational reason to  
10 suppose the claims are literally accurate. We already know from  
11 cognitive neuroscience that the brain constructs many internal models  
12 — bundles of information that represent items in the real world  
13 (Johnson-Laird, 1983; Holmes and Spence, 2004; Graziano, 2013).  
14 These models, whenever they have been studied in detail, are always  
15 simplified. They are quick-and-dirty descriptions, useful if not entirely  
16 accurate. The question in front of us is not: how does the brain  
17 generate a non-physical essence? Rather, we should ask: what set of  
18 information in the brain is the basis for our claim to have conscious  
19 experience, and what adaptive function does that information serve?  
20 AST does not explain how a brain generates a subjective experience.  
21 It explains how a machine makes claims about itself, and how the  
22 information on which those claims are based may have a cognitive,  
23 functional use.

24 Chalmers has written an insightful article, outlining what he has  
25 termed the hard problem and the meta-problem (Chalmers, 2018). One  
26 way to frame the hard problem is that consciousness is a private  
27 experience whose existence cannot be assessed from the outside.  
28 Because it cannot be physically measured, it cannot be scientifically  
29 studied. The meta-problem, in contrast, is the question of why we  
30 think we have a hard problem. Part of Chalmers' discussion focuses  
31 on an approach to consciousness called illusionism (Frankish, 2016).  
32 In that approach, consciousness does not exist as such — it is illusory.  
33 One of the earliest and most influential illusionist accounts is  
34 Dennett's idea of the user illusion (Dennett, 1992). Illusionism could  
35 be considered a proposed approach to the meta-problem — it suggests  
36 that we think we have a hard problem of consciousness because we  
37 are misinformed by an illusion.

38 AST specifically addresses Chalmers' meta-problem, because it  
39 addresses how a biological machine claims to have a hard problem.  
40 Yet in Chalmers' article, one senses his uneasiness over how to inter-  
41 pret AST. For example, he puzzles over the question: in AST, what

1 exactly is awareness? Is it an attention schema, or is it supposed to be  
 2 an abstraction to which an attention schema refers? As I spell out in  
 3 chapter 3 of my book *Consciousness and the Social Brain* (2013),  
 4 AST does not easily pin down what, exactly, awareness itself is. The  
 5 reason for the ambiguity, I believe, is that AST is fundamentally not a  
 6 philosophical theory. It is an engineering theory. It explains the  
 7 performance of a machine — it explains how a machine claims to  
 8 have consciousness. It could be viewed as an illusionist theory, and is  
 9 especially close to Dennett’s account. Yet it may not perfectly fit into  
 10 the illusionist category either — or at least it may provide a different  
 11 emphasis. Illusionism seems to ask: how does the brain generate, if  
 12 not an actual conscious experience, at least an illusory semblance of  
 13 one? That framing focuses on how the brain generates something, and  
 14 on consciousness as a distinct item of interest whose real or illusory  
 15 nature can be debated. But in AST there is no meaningful answer to  
 16 the question. Instead, the theory addresses how a machine makes  
 17 claims, not how a machine generates experiences or illusions. We can  
 18 understand how a car drives and a bird flies, from an engineering  
 19 perspective. We should be able to understand, mechanistically, how a  
 20 brain makes claims.

## 21 **2. Model-Based Knowledge**

22 To describe consciousness as the brain making claims, I acknowledge,  
 23 sounds at first too reductive. But the crux of the argument lies in the  
 24 information sets on which those claims are based. AST depends on  
 25 model-based knowledge, as distinct from superficial knowledge. To  
 26 explain what I mean, I will use an example that I have used in other  
 27 recent accounts (Graziano, 2019).

28 Suppose a child plays at make-believe. She barks, crawls on all  
 29 fours, and says, ‘I’m a puppy!’ Something in her brain contains the  
 30 information that puppies bark and walk on all fours. Her brain has also  
 31 constructed the proposition ‘I’m a puppy!’ or else she would not be  
 32 able to make the claim. And yet that information exists in a larger  
 33 context. Her brain contains a net of information including ‘I’m not  
 34 really a puppy’, ‘I’m making it up’, ‘I’m a little girl’, and so on. Some  
 35 of that information is present at a cognitive and linguistic level. Much  
 36 of it is at a deeper, sensory or perceptual level. Her body schema is  
 37 constructed automatically, beneath higher cognition, and describes the  
 38 physical layout of a human body, not a puppy body. She sees her  
 39 human hands in front of her, and the representations constructed in her

1 visual system confirm her human identity. She remembers eating  
2 breakfast with a spoon, going to school, reading a book — all human  
3 activities. The claim ‘I’m a puppy’ is superficial knowledge that is  
4 inconsistent with her deepest internal models.

5 But suppose I have the science fiction tools to manipulate the  
6 information in her brain. I alter her body schema to reflect the body of  
7 a puppy. I alter the information in her visual system and her memory  
8 to make it consistent with the puppy proposition. I remove the specific  
9 cognitive information that says ‘I made that up to play a game’. I  
10 switch the information that says ‘I’m certain this is not true’ to its  
11 opposite. How would she know that she is not a puppy? Her brain is  
12 captive to the information it contains. Tautologically, it knows what it  
13 knows. She would no longer think of her puppy identity as a hypo-  
14 theoretical. She would take it as a literal truth. There would be no reason  
15 for her to think otherwise. One might say that she now believes,  
16 intuitively, that she is a puppy; and here, to clarify the terminology, by  
17 ‘believing something intuitively’ I mean that her cognition is informed  
18 by deeper, automatically constructed, internal models. The belief, at  
19 the cognitive level, derives from the deeper internal models over  
20 which she has no cognitive control.

21 You could tell her, ‘But you understand English. Puppies can’t do  
22 that. Don’t you think that suggests you’ve mistaken your identity?’ If  
23 she is intellectually precocious, she might realize the logic of your  
24 argument. That new information, however, will be at a superficial,  
25 cognitive level. It will conflict with her deeper internal models. Like  
26 so many people, she will be in a position of believing one truth about  
27 herself intuitively, while entertaining a different truth intellectually.

28 Just so, I might be able to convince you intellectually that your  
29 claim to consciousness has its basis in an information set — an  
30 attention schema, as I’ll explain in the next section. But intuitively,  
31 you still believe a different truth about yourself. When you rely on  
32 introspection — when your cognition accesses deeper internal models  
33 — they provide you with a different story. They inform you  
34 (incorrectly) that your consciousness is not just information or compu-  
35 tation — it has a ‘what it feels like’ component, an ethereal essence  
36 dwelling inside you. Even if I have convinced you of my argument,  
37 you will find yourself conflicted, with superficial, intellectual knowl-  
38 edge pointing you towards one understanding and deeper, internal  
39 models, over which you have no cognitive control, anchoring you to a  
40 different understanding.

1

### 3. The Attention Schema

2 In this short piece, I will not give a complete account of AST or the  
3 supporting lines of evidence. I refer readers to previous publications  
4 (e.g. Graziano and Kastner, 2011; Graziano, 2013; Webb and  
5 Graziano, 2015). Instead, here, I briefly summarize the core concept.

6 Logically, we claim to have subjective experience for the same  
7 reason we make any claim — because the brain has constructed the  
8 requisite information on which the claim is based. Suppose a person  
9 looks at a red apple and reports having a subjective experience of red.  
10 It is not enough for the brain to construct colour information, which  
11 would allow the person to make the limited claim ‘The apple is red’.  
12 We know, for example, that people who suffer from blindsight  
13 (Cowey, 2010) can process visual information and make claims about  
14 visual features, without reporting any conscious visual experience. To  
15 report a conscious experience, the brain must also construct the  
16 information on which it bases the claim ‘I have something extra, a  
17 non-physical subjective experience, associated with the redness’.

18 The brain constructs descriptive sets of information because they act  
19 as useful models for real items in the world. One question facing us,  
20 therefore, is: what is the physically real item that is modelled by this  
21 particular information set, on which the claim of conscious experience  
22 is based?

23 For example, your brain constructs a set of information that is,  
24 moment by moment, correlated with the configuration of your right  
25 arm. That information set is called an arm schema, a part of the body  
26 schema (Graziano and Botvinick, 2002; Holmes and Spence, 2004; de  
27 Vignemont, 2018). It is the basis on which you can close your eyes  
28 and report on the presence and state of your arm. That internal model  
29 usually covaries with the physical arm, although the two can be  
30 dissociated. Like all internal models in the brain, the model of the arm  
31 is a detail-poor simplification, and can sometimes make errors and  
32 become misaligned. It *usually* describes the overall state of the arm.  
33 One could say that the fact that this particular set of signals in the  
34 brain co-varies with the state of the arm, by definition, makes it an  
35 arm schema. The close tracking of the arm is what makes it informa-  
36 tive about the arm.

37 Can we find any physically real, objectively measurable item that  
38 co-varies with people’s report of conscious experience? Yes. This  
39 question has a straightforward answer known in psychology and  
40 neuroscience for decades. The report of conscious experience tends to

1 co-vary with attention (e.g. Posner, 1994; Mack and Rock, 1998;  
2 Simons and Chabris, 1999; Cohen *et al.*, 2012). If a person reports  
3 being conscious of X, she is typically also attending to X. Attention  
4 and awareness can sometimes be separated. At least, attention without  
5 awareness has been demonstrated, though awareness without attention  
6 has not yet been convincingly shown (e.g. Kentridge, Heywood and  
7 Weiskrantz, 1999; Tsushima, Sasaki and Watanabe, 2006; Webb,  
8 Kean and Graziano, 2016). Separating the two depends on pushing the  
9 system to extremes, either through brain damage or laboratory con-  
10 ditions in which visual stimuli are degraded and presented at detection  
11 threshold. Most of the time, however, awareness closely tracks  
12 attention. (Indeed it seems to be easier to separate the arm from the  
13 arm schema than attention from awareness, at least in my experience  
14 having experimentally studied both topics.)

15 One might ask: is attention too narrow a phenomenon to cover sub-  
16 jective consciousness? Surely we are conscious of much more than we  
17 put at the focus of our attention. But objectively speaking, in decades  
18 of work on subjective awareness and attention, this intuition is not  
19 correct (e.g. Posner, 1994; Mack and Rock, 1998; Simons and  
20 Chabris, 1999; Cohen *et al.*, 2012). Awareness and attention co-vary  
21 most of the time. The confusion arises when people use a colloquial  
22 definition of attention, rather than a scientific one. In a typical collo-  
23 quial definition, attention is a limited, central focus within the larger  
24 field of consciousness. In contrast, in neuroscience and psychology,  
25 attention is a process in the brain, primarily in the cerebral cortex,  
26 whereby a representation (such as a visual representation of an apple)  
27 has its signals enhanced, competing representations have their signals  
28 reduced, and the enhanced signals have a correspondingly greater  
29 impact on systems around the brain (Desimone and Duncan, 1995;  
30 Beck and Kastner, 2009). That enhancement can occur either due to  
31 greater external salience (bottom-up attention) or due to internal  
32 modulation (top-down control). Attention is not limited to one central  
33 object; it can be directed away from the fovea, for example, and it can  
34 be spread and divided. If you think that you are aware of something  
35 outside of your attention — that you are attending only to A while  
36 also aware of B, C, and D — that intuition is not correct; or at least,  
37 you are drawing on a colloquial definition of attention. By the  
38 scientific definition, you are probably attending to all of these items to  
39 some degree. Consciousness almost always co-varies with attention. It  
40 therefore effectively serves as a model of attention.

1 Chalmers suggests that linking consciousness to an attention schema  
2 is overly specific. Perhaps the brain constructs a general ‘representa-  
3 tion schema’ which tells us what it means to represent information and  
4 gives rise to our claims about consciousness. But this suggestion  
5 stems from a misunderstanding of the theory. The report of conscious  
6 experience does not correlate with all representations in the brain. It  
7 correlates specifically with attention. Just so, the internal model of my  
8 arm is not a general ‘moving object schema’. It is specifically an arm  
9 schema, because it tracks the state of my arm. Moreover, the func-  
10 tional use of an arm schema is to monitor, predict, and help control  
11 your arm; and the proposed functional use of an attention schema is to  
12 monitor, predict, and help control attention.

13 Suppose we were to design an attention schema from scratch. Our  
14 goal is to construct a useful information set descriptive of attention.  
15 For comparison, the arm schema contains stable information such as  
16 size, shape, jointed structure, and weight, as well as changing informa-  
17 tion such as how the arm is moving at the moment. Just so, the  
18 attention schema might describe both stable and changing properties  
19 of attention. Imagine a rich, textbook-style, scientific description of  
20 attention, including the details of the physical mechanisms present in  
21 the brain — and then imagine stripping from that description every-  
22 thing unnecessary for the brain to be informed about. We strip away  
23 information about neurons, synapses, inhibition and excitation — the  
24 physical truth of attention. We strip away information about bottom-  
25 up and top-down pathways, about fronto-parietal networks, about the  
26 thalamus and about the superior colliculus. We strip away information  
27 about the technical distinctions between exogenous and endogenous,  
28 engage and disengage, overt and covert, spatial and feature. We are  
29 left with a detail-poor description of attention as an amorphous ‘thing’  
30 inside of me, a mental stuff that can grasp hold of objects in an  
31 abstract sense. The ‘thing’ can grasp hold of external objects like an  
32 apple, or internal objects like the thought that  $2 + 2 = 4$ . The ‘thing’  
33 has special powers such that, when it grasps hold of object X, it causes  
34 me to understand the details and the deeper meaning of X; it causes X  
35 to become vivid to me; it empowers me to choose to react to X, and to  
36 remember it for later. This stripped-down description of attention  
37 contains no information about the physical properties of the ‘thing’  
38 inside me. As far as one can tell from the attention schema, that  
39 ‘thing’ lacks physicality.

40 My argument here is that if a brain uses the mechanism of attention,  
41 and if it constructs a simplified internal model of it, and if it makes

1 claims about itself on the basis of the information in that attention  
2 schema, then it ought to claim to have a subjective, non-physical,  
3 mental grasp, or experience, of objects. In this way, AST explains  
4 how a machine claims to have consciousness — without having to  
5 explain what consciousness itself is.

#### 6 **4. The Non-physical Essence**

7 The philosopher François Kammerer asked an insightful question  
8 (Kammerer, 2016; 2018). Suppose AST is correct. The brain con-  
9 structs an attention schema which represents general properties of  
10 attention, such as our ability to focus on and process information in  
11 depth. At the same time, it leaves out any depiction of the physical or  
12 mechanistic properties of attention. It does not specify that attention  
13 *lacks* a physical substance — it is merely silent on the topic. It is  
14 uninformative on the details of neurons and synapses. If our claims  
15 about consciousness derive from that internal model, then why do  
16 people typically make such a strong claim that consciousness is an  
17 ethereal essence, something inside of us that specifically *lacks*  
18 physical substance? Why do we not, instead, have an intuition of  
19 consciousness as an entity whose physical attributes — weight, size,  
20 hardness — are simply not yet known?

21 The answer may lie partly in a subtle distinction. I suggest that we  
22 do not generally understand consciousness as a thing whose physical  
23 dimensions are undetermined. Instead, we intuitively understand con-  
24 sciousness as something for which physical dimensions are *irrelevant*.

25 Imagine someone taps you on the shoulder. The touch activates skin  
26 receptors, and neuronal fibres transmit that information to the brain.  
27 Ultimately, your brain constructs a specific kind of internal model, a  
28 tactile model, a packet of information that describes that particular  
29 touch. The model contains information about the location of the touch,  
30 the intensity at onset, the pressure, the duration, the smooth or plush  
31 texture of a fingertip. It is a rich sensory representation. But it contains  
32 no information about taste. A touch on the shoulder does not come  
33 with a salty taste. I do not mean that a touch is bland and needs salt —  
34 no, it does not lie *anywhere* on any taste dimension. It does not  
35 occupy the same information space. Now that I have mentioned the  
36 possibility, you can consider it in a superficial, cognitive sense, but  
37 you cannot alter the deeper, internal model. Touch perception is an  
38 inborn process and is not open to cognitive modification.



1 If you could insert electrodes into a person's brain and read the  
2 information encoded in the tactile system, the perceptual model for a  
3 touch would presumably not contain the information 'And by the way,  
4 no taste is present'. It does not need the explicit negation. It is simply  
5 silent on taste. We do not intuitively understand touch to be something  
6 for which taste has been minimized; or something that might have a  
7 taste, but we just don't know yet what the taste is. Instead, we under-  
8 stand touch to be something for which taste is *irrelevant*.

9 I argue that the attention schema acts the same way. It depicts  
10 general properties of attention, but not physical, mechanistic  
11 properties. Based on that internal model, we intuitively believe in an  
12 inner mental experience that takes possession of information and  
13 drives action, the way attention does, but that has no specific relation-  
14 ship to physicality. Physicality is irrelevant to it. That mental essence  
15 is not physically graspable, smooth, textured, rough, bumpy, heavy,  
16 light, smelly, green, pointy — it does not lie anywhere on those  
17 physical dimensions, any more than a touch exists on the salty  
18 dimension.

19 And yet, in AST, the attention schema depicts at least one physical  
20 property. It depicts attention as having a physical location roughly  
21 inside us (see my prior accounts of the importance of localization in a  
22 model of attention: Graziano and Kastner, 2011; Graziano, 2013).  
23 Based on the information within that internal model, we should have  
24 an intuition about a mental essence that overlaps the physical world, in  
25 that you can point to a location and say 'it lives roughly here'. It is  
26 like a ghost, inhabiting physical space even as it lacks any relationship  
27 to other physical attributes. It has its own special power — to make us  
28 know and react. In this theory, the ghost in the machine, the con-  
29 sciousness inside us, is a topic of discussion among us only because  
30 our intuitions are informed by an attention schema, with its incom-  
31 plete account of attention.

32 And so we come back to the hard problem and the meta-problem. In  
33 my proposed explanation, the belief in a hard problem derives from  
34 intuitions that come bubbling up from a deep, subsurface model, the  
35 attention schema. AST is a meta answer that explains why people  
36 believe in a hard problem in the first place.

## 37 References

38 Beck, D.M. & Kastner, S. (2009) Top-down and bottom-up mechanisms in biasing  
39 competition in the human brain, *Vision Research*, **49**, pp. 1154–1165.

- 1 Chalmers, D.J. (2018) The meta-problem of consciousness, *Journal of Consciousness Studies*, **25** (9–10), pp. 6–61.
- 2
- 3 Cohen, M.A., Cavanagh, P., Chun, M.M. & Nakayama, K. (2012) The attentional requirements of consciousness, *Trends in Cognitive Sciences*, **16**, pp. 411–417.
- 4
- 5 Cowey, A. (2010) The blindsight saga, *Experimental Brain Research*, **200**, pp. 3–24.
- 6
- 7 de Vignemont, F. (2018) *Mind and Body*, New York: Oxford University Press.
- 8 Dennett, D.C. (1992) *Consciousness Explained*, New York: Little-Brown.
- 9 Desimone, R. & Duncan, J. (1995) Neural mechanisms of selective visual attention, *Annual Review of Neuroscience*, **18**, pp. 193–222.
- 10
- 11 Frankish, K. (2016) Illusionism as a theory of consciousness, *Journal of Consciousness Studies*, **23** (11–12), pp. 11–39. Reprinted in Frankish, K. (ed.) (2017) *Illusionism as a Theory of Consciousness*, Exeter: Imprint Academic.
- 12
- 13
- 14 Graziano, M.S.A. (2013) *Consciousness and the Social Brain*, Oxford: Oxford University Press.
- 15
- 16 Graziano, M.S.A. (2019) *Rethinking Consciousness*, New York: W.W. Norton.
- 17
- 18 Graziano, M.S.A. & Botvinick, M.M. (2002) How the brain represents the body: Insights from neurophysiology and psychology, in Prinz, W. & Hommel, B. (eds.) *Common Mechanisms in Perception and Action: Attention and Performance XIX*, pp. 136–157, Oxford: Oxford University Press.
- 19
- 20
- 21 Graziano, M.S.A. & Kastner, S. (2011) Human consciousness and its relationship to social neuroscience: A novel hypothesis, *Cognitive Neuroscience*, **2**, pp. 98–113.
- 22
- 23
- 24 Holmes, N.P. & Spence, C. (2004) The body schema and the multisensory representation(s) of peripersonal space, *Cognitive Processing*, **5**, pp. 94–105.
- 25
- 26 Johnson-Laird, P. (1983) *Mental Models*, New York: Lawrence Erlbaum.
- 27
- 28 Kammerer, F. (2016) The hardest aspect of the illusion problem — and how to solve it, *Journal of Consciousness Studies*, **23** (11–12), pp. 124–139. Reprinted in Frankish, K. (ed.) (2017) *Illusionism as a Theory of Consciousness*, Exeter: Imprint Academic.
- 29
- 30
- 31 Kammerer, F. (2018) Can you believe it? Illusionism and the illusion meta-problem, *Philosophical Psychology*, **31**, pp. 44–67.
- 32
- 33 Kentridge, R.W., Heywood, C.A. & Weiskrantz, L. (1999) Attention without awareness in blindsight, *Proceedings of the Royal Society B: Biological Sciences*, **266**, pp. 1805–1811.
- 34
- 35
- 36 Mack, A. & Rock, I. (1998) *Inattentional Blindness*, Cambridge, MA: MIT Press.
- 37
- 38 Posner, M.I. (1994) Attention: The mechanisms of consciousness, *Proceedings of the National Academy of Sciences USA*, **91**, pp. 7398–7403.
- 39
- 40 Simons, D.J. & Chabris, C.F. (1999) Gorillas in our midst: Sustained inattention blindness for dynamic events, *Perception*, **28**, pp. 1059–1074.
- 41
- 42 Tsushima, Y., Sasaki, Y. & Watanabe, T. (2006) Greater disruption due to failure of inhibitory control on an ambiguous distractor, *Science*, **314**, pp. 1786–1788.
- 43
- 44 Webb, T.W. & Graziano, M.S.A. (2015) The attention schema theory: A mechanistic account of subjective awareness, *Frontiers in Psychology*, **6**, art. 500.
- 45
- 46
- 47 Webb, T.W., Kean, H.H. & Graziano, M.S.A. (2016) Effects of awareness on the control of attention, *Journal of Cognitive Neuroscience*, **28**, pp. 842–851.